

Delay Analysis for a Discretionary-Priority Packet-Switching System

Sung Jo Hong (洪誠條)* and Hideaki Takagi (高木英明)**

**Institute of Information Sciences and Electronics, University of Tsukuba, Japan*

***Institute of Socio-Economic Planning, University of Tsukuba, Japan*

Abstract

We consider a priority-based packet-switching system with three phases of the packet transmission time. Each packet belongs to one of several priority classes, and the packets of each class arrive at a switch in a Poisson process. The switch transmits queued packets on a priority basis with three phases of preemption mechanism. Namely, the transmission time of each packet consists of a preemptive-repeat part for the header, a preemptive-resume part for the information field, and a nonpreemptive part for the trailer.

By an exact analysis of the associated queueing model, we obtain the Laplace-Stieltjes transform of the distribution function for the delay, i.e., the time from arrival to transmission completion, of a packet for each class. We derive a set of equations that calculates the mean response time for each class recursively. Based on this result, we plot the numerical values of the mean response times for several parameter settings. The probability generating function and the mean for the number of packets of each class present in the system at an arbitrary time are also given.

1. Introduction

Priority-based transmission can be used in packet switching of mixed traffic in order to handle properly different characteristics of service requirements for several classes of packets, including controlling commands and responses, real-time data, file transfer, and so on. In fact, in the network layer of Systems Network Architecture (SNA), messages are placed on the appropriate transmission group outbound queue in a priority sequence determined by the transmission priority field and network priority indicator in the format identifier type 4 (FID4) transmission header, where three levels of priority can be specified [1]. In the IEEE 802.5 standard for the token passing ring for local area networks, eight levels of priority can be specified in the starting control field of the frame format, and a frame of the highest priority in the network is given the right for transmission [10]. The fiber distributed data interface (FDDI) and the distributed queue dual bus (DQDB) for metropolitan area networks also implement similar priority mechanisms. These are examples of the nonpreemptive priority service discipline, which means that the service of a packet is continued until completion even if other packets with higher priority arrive during its service.

Recently, Cho and Un [4] (also Cho [3]) proposed and analyzed a queueing system with a combined preemptive/nonpreemptive priority discipline as a model of message transmission strategies for communication systems and job scheduling in computer systems. Their idea, which originates in the *discretionary priority queueing* of Avi-Itzhak et al. [2] and Jaiswal [7, Chapter VI], is to split the packet transmission time into two phases such that the first phase is served according to the preemptive resume, repeat-identical or repeat-different discipline while the second phase is served in the nonpreemptive manner. Earlier, Komatsu [9] studied a two-priority queueing system with three phases of preemption mechanism. Namely, the transmission time of each packet consists of a preemptive repeat-different part for the header, a preemptive resume part for the information field, and a nonpreemptive part for the trailer. According to Komatsu, the motivation for this phasing is as follows. Since the header of a packet includes control information for the packets, it does not make sense to retransmit the preempted header from the point of preemption; the whole header must be retransmitted. If the transmission is preempted during the information field, it can be resumed from the point of preemption, because the control information has already been received. Finally, it

would not be efficient to allow preemption of the transmission during the trailer, because the whole transmission has almost been completed.

The present paper unifies and extends the works of Cho and Un [4] and Komatsu [9] by considering a queueing system for packets of several priority classes such that the transmission time of each packet is split into a preemptive repeat-identical, a preemptive resume part, and a nonpreemptive part, where the joint probability distribution for the lengths of the three parts is generally given. Our model is more general than that of Cho and Un who deal with a system of priority packets with two preemption phases and special splitting schemes, and that of Komatsu who considers a system with only two priority levels and independent lengths of the three phases. All these studies assume that the packets of each priority class arrive in a Poisson process and that the transmission time of each packet has general distribution.

2. Model and Notation

We consider a multiclass priority queueing system. Let there be P classes of packets indexed as $1, 2, \dots, P$. Packets of class p arrive in a Poisson process at rate λ_p . We assume that the classes of packets are priority classes such that class p has priority over class q if $p < q$. The service time of a packet of class p consists of three phases, namely, the *preemptive repeat* (RP) phase, the *preemptive resume* (RS) phase, and the *nonpreemptive* (NP) phase in this order. Packets are preferentially served by a single server in the order of priority, and for each priority in the order of arrival. If a packet of high priority arrives when a packet of lower priority is being served in its RP phase, the server interrupts the current service and immediately starts to serve the packet of high priority. The preempted service for the packet of lower priority is commenced again from the beginning, following one of the two types of repeat discipline: if the service time to be repeated is the same amount of service as the one interrupted, the discipline is called *preemptive repeat-identical*; if it is newly sampled from the given distribution for the preempted class, this is called *preemptive repeat-different*. We will consider the preemptive repeat-identical discipline only. If a packet of high priority arrives when a packet of lower priority is being served in its RS phase, the server interrupts the current service and immediately starts to serve the packet of high priority. The preempted service for the packet of lower priority is commenced again from the point where it was interrupted. If a packet of high priority arrives when a packet of lower priority is being served in its NP phase, the server never interrupts the current service. Thus the high priority packet waits in the queue until the end of the current service.

In this paper, we obtain the distributions for the *residence time* R_p , the *waiting time* W_p and the *response time* T_p of a packet of class p , and the distribution for the number L_p of packets of class p present in the system at an arbitrary time. The residence time R_p is the interval from the start to the completion of its service, including periods during which the service is preempted. The waiting time W_p is the interval from the arrival to the service start of a packet of class p . The response time T_p is the interval from the arrival to the entire service completion of a packet of class p . T_p is the sum of W_p and R_p , which are independent of each other. Note the relation

$$T_p = W_p + R_p, \quad p = 1, \dots, P, \quad (2.1)$$

where W_p and R_p are independent. As preliminaries, we analyze the *completion time* C_p for a packet of class p , which is the interval from the start of its service to the first moment after the service completion at which no packets of class 1 through $p - 1$ are present in the system, and the length Θ_p^+ of a *busy period* generated by packets of class 1 through p .

Let x_p^{RP} , x_p^{RS} , and x_p^{NP} be the random variables that denote the lengths of the RP phase, the RS phase, and the NP phase of the service time for a packet of class p , respectively. Then the entire service time of a packet of class p is given by $x_p = x_p^{RP} + x_p^{RS} + x_p^{NP}$. We assume that x_p^{RP} , x_p^{RS} , and x_p^{NP} can be dependent. The marginal distribution functions (DFs) for x_p^{RP} , x_p^{RS} , and x_p^{NP} are denoted by $B_p^{RP}(x)$, $B_p^{RS}(x)$, and $B_p^{NP}(x)$, and the corresponding Laplace-Stieltjes transforms

(LSTs) of the DF by $B_p^{*,RP}(s)$, $B_p^{*,RS}(s)$, and $B_p^{*,NP}(s)$, respectively. Let \bar{x}_p be the *gross service time* that the server attends to a packet of class p . The server utilization of packets of class p is given by

$$\rho_p = \lambda_p E[\bar{x}_p], \quad (2.2)$$

where $E[\bar{x}_p]$ is the expected value of \bar{x}_p . The aggregate arrival rate of packets of class 1 through p and their server utilization are given by

$$\lambda_p^+ = \sum_{k=1}^p \lambda_k, \quad \rho_p^+ = \sum_{k=1}^p \rho_k \quad (2.3)$$

Similarly, the aggregate arrival rate of packets of class $p+1$ through P and their server utilization are given by

$$\lambda_p^- = \sum_{k=p+1}^P \lambda_k, \quad \rho_p^- = \sum_{k=p+1}^P \rho_k \quad (2.4)$$

The total server utilization is given by $\rho = \sum_{k=1}^P \rho_k$. Throughout the paper we assume that the system is not saturated in the steady state, i.e., $\rho < 1$.

3. Completion time

Let us consider the completion time C_p for a packet of class p . The completion time C_p can be decomposed into two periods: (i) the period C_p^{RP} from the moment at which a packet of class p first receives its service until the time when it completes the RP phase of its service, and (ii) the period C_p^{RS+NP} from the moment at which a packet of class p first receives the RS phase of its service to the first moment after its service completion at which there are no packets of class $1, 2, \dots, p-1$ in the system.

We first consider the period C_p^{RS+NP} . The period C_p^{RS+NP} may be regarded as a delay cycle with the initial delay $x_p^{RS+NP} \equiv x_p^{RS} + x_p^{NP}$ generated by packets of class $1, 2, \dots, p-1$. Therefore we have

$$E[e^{-sC_p^{RS+NP}} | x_p^{RS+NP}] = e^{-(s+\lambda_{p-1}^+ - \lambda_{p-1}^+ \Theta_{p-1}^+(s))x_p^{RS+NP}} \quad (3.1)$$

We next consider the period C_p^{RP} . In the case of preemptive repeat-identical discipline, let $x_p^{RP}(n)$ be the service time futilely expended by the n^{th} preemption. Let $\Theta_{p-1}^+(n)$ be the duration of the n^{th} preemption. Note that $\Theta_{p-1}^+(n)$ is the busy period of packets of class $1, 2, \dots, p-1$. Then the period C_p^{RP} can be written as follows:

$$C_p^{RP} = x_p^{RP} + \sum_{n=1}^N x_p^{RP}(n) + \sum_{n=1}^N \Theta_{p-1}^+(n) \quad (3.2)$$

Given x_p^{RP} , the distribution of the number of preemptions is given by

$$P\{N = n | x_p^{RP}\} = (1 - e^{-\lambda_{p-1}^+ x_p^{RP}})^n e^{-\lambda_{p-1}^+ x_p^{RP}}, \quad n = 0, 1, 2, \dots \quad (3.3)$$

Given that the phase x_p^{RP} is preempted, the distribution of the expended service time $x_p^{RP}(n)$ is given by

$$P\{x_p^{RP}(n) \leq x | x_p^{RP} \text{ is preempted}\} = \frac{1 - e^{-\lambda_{p-1}^+ x}}{1 - e^{-\lambda_{p-1}^+ x_p^{RP}}}, \quad 0 \leq x \leq x_p^{RP} \quad (3.4)$$

which yields

$$E[e^{-sx_p^{RP}(n)} | x_p^{RP} \text{ is preempted}] = \frac{1}{1 - e^{-\lambda_{p-1}^+ x_p^{RP}}} \int_0^{x_p^{RP}} e^{-sx} \lambda_{p-1}^+ e^{-\lambda_{p-1}^+ x} dx \quad (3.5)$$

$$= \frac{\lambda_{p-1}^+ (1 - e^{-(s+\lambda_{p-1}^+) x_p^{RP}})}{(s + \lambda_{p-1}^+) (1 - e^{-\lambda_{p-1}^+ x_p^{RP}})} \quad (3.6)$$

Using (3.6) in (3.2), we get

$$E[e^{-sC_p^{RP}} | x_p^{RP}, N] = e^{-sx_p^{RP}} \left(\frac{\lambda_{p-1}^+}{s + \lambda_{p-1}^+} \right)^N \left(\frac{1 - e^{-(s+\lambda_{p-1}^+) x_p^{RP}}}{1 - e^{-\lambda_{p-1}^+ x_p^{RP}}} \right)^N (\Theta_{p-1}^+(s))^N \quad (3.7)$$

where $\Theta_{p-1}^+(s)$ is the LST of the DF for $\Theta_{p-1}^+(n)$, which is independent of n . Removing the condition on N by using (3.3), we get

$$E[e^{-sR_p^{RP}} | x_p^{RP}] = \frac{(s + \lambda_{p-1}^+) e^{-(s+\lambda_{p-1}^+) x_p^{RP}}}{s + \lambda_{p-1}^+ - \lambda_{p-1}^+ \Theta_{p-1}^+(s) [1 - e^{-(s+\lambda_{p-1}^+) x_p^{RP}}]} \quad (3.8)$$

Thus, from (3.1) and (3.8), if the RP phase for a packet of class p is served according to the preemptive repeat-identical discipline, we obtain the LST $C_p^*(s)$ for the completion time C_p as

$$C_p^*(s) = \int_0^\infty \int_0^\infty \frac{(s + \lambda_{p-1}^+) e^{-(s+\lambda_{p-1}^+) x} e^{-(s+\lambda_{p-1}^+ - \lambda_{p-1}^+ \Theta_{p-1}^+(s)) y}}{s + \lambda_{p-1}^+ - \lambda_{p-1}^+ \Theta_{p-1}^+(s) [1 - e^{-(s+\lambda_{p-1}^+) x}]} dB_p^{RP,RS+NP}(x, y) \quad (3.9)$$

where $B_p^{RP,RS+NP}(x, y)$ is the joint DF for x_p^{RP} and x_p^{RS+NP} . From (3.9), we get

$$E[C_p] = \left(\frac{1}{\lambda_{p-1}^+} + E[\Theta_{p-1}^+] \right) \left(E[e^{\lambda_{p-1}^+ x_p^{RP}}] + \lambda_{p-1}^+ E[x_p^{RS+NP}] - 1 \right) \quad (3.10)$$

$$\begin{aligned} E[C_p^2] &= 2 \left(\frac{1}{\lambda_{p-1}^+} + E[\Theta_{p-1}^+] \right)^2 E[(e^{\lambda_{p-1}^+ x_p^{RP}} - 1)^2] \\ &\quad + \left(\frac{2}{(\lambda_{p-1}^+)^2} + \frac{2E[\Theta_{p-1}^+]}{\lambda_{p-1}^+} + E[(\Theta_{p-1}^+)^2] \right) \left(E[e^{\lambda_{p-1}^+ x_p^{RP}}] - 1 \right) \\ &\quad - 2 \left(\frac{1}{\lambda_{p-1}^+} + E[\Theta_{p-1}^+] \right) E[x_p^{RP} e^{\lambda_{p-1}^+ x_p^{RP}}] \\ &\quad + (1 + \lambda_{p-1}^+ E[\Theta_{p-1}^+])^2 \\ &\quad \times \left(E[(x_p^{RS+NP})^2] + \frac{2E[x_p^{RS+NP} e^{\lambda_{p-1}^+ x_p^{RP}}]}{\lambda_{p-1}^+} - \frac{2E[x_p^{RS+NP}]}{\lambda_{p-1}^+} \right) \\ &\quad + \lambda_{p-1}^+ E[(\Theta_{p-1}^+)^2] E[x_p^{RS+NP}]. \end{aligned} \quad (3.11)$$

4. Busy period

Let us first derive the relation between completion times and busy periods. (See Jaiswal [7, Section IV. 7-2] and Takagi [11, Section 3.4] for the same treatment for systems with only preemptive repeat priority discipline.)

We first consider $\Theta_p^+(s)$, the LST of the DF for the length Θ_p^+ of a busy period generated by packets of class $1, 2, \dots, p$. If the packet that arrives first in an idle period is of class p , which occurs

with probability λ_p/λ_p^+ , the busy period Θ_p^+ is equal to a busy period for packets of class p that have the completion time C_p as the service time. Let Θ_p be the length of a busy period generated by packets of class p that have the "service time" C_p . Then the LST $\Theta_p^*(s)$ of the DF for Θ_p satisfies the equation

$$\Theta_p^*(s) = C_p^*(s + \lambda_p - \lambda_p \Theta_p^*(s)). \quad (4.1)$$

If the first arriving packet is of class $1, 2, \dots, p-1$, which occurs with probability $\lambda_{p-1}^+/\lambda_p^+$, then the busy period Θ_p^+ is equal to the delay cycle with initial delay Θ_{p-1}^+ generated by packets of class p . Thus we have

$$\Theta_p^+(s) = \frac{\lambda_p}{\lambda_p^+} \Theta_p^*(s) + \frac{\lambda_{p-1}^+}{\lambda_p^+} \Theta_{p-1}^+(s + \lambda_p - \lambda_p \Theta_p^*(s)). \quad (4.2)$$

From (4.1) and (4.2) we get the recursive relations

$$E[\Theta_p^+] = \frac{\lambda_p E[C_p] + \lambda_{p-1}^+ E[\Theta_{p-1}^+]}{\lambda_p^+ (1 - \lambda_p E[C_p])} \quad (4.3)$$

$$E[(\Theta_p^+)^2] = \frac{\lambda_p (1 + \lambda_{p-1}^+ E[\Theta_{p-1}^+]) E[C_p^2]}{\lambda_p^+ (1 - \lambda_p E[C_p])^3} + \frac{\lambda_{p-1}^+ E[(\Theta_{p-1}^+)^2]}{\lambda_p^+ (1 - \lambda_p E[C_p])^2} \quad (4.4)$$

5. Residence time

We derive the LSTs of the DF for the residence time and the gross service time for a packet of class p . The residence time R_p for a packet of class p consists of the service time x_p^{NP} of its NP phase and the following two periods: (i) the period R_p^{RP} from the moment at which a packet of class p first receives its service until the time when it completes the RP phase of its service, and (ii) the period R_p^{RS} from the moment at which a packet of class p first receives the RS phase of its service until the time when it completes the RS phase.

We first consider the period R_p^{RS} . Suppose that N preemptions occur because of the arrivals of packets of class $1, 2, \dots, p-1$ during the RS phase. The period R_p^{RS} for a packet of class p can then be written as

$$R_p^{RS} = x_p^{RS} + \sum_{n=1}^N \Theta_{p-1}^+(n) \quad (5.1)$$

Given x_p^{RS} and N , we get

$$E[e^{-sR_p^{RS}} | x_p^{RS}, N] = e^{-sx_p^{RS}} [\Theta_{p-1}^+(s)]^N \quad (5.2)$$

Given x_p^{RS} , the distribution of the number of preemptions is given by

$$P\{N = n | x_p^{RS}\} = \frac{(\lambda_{p-1}^+ x_p^{RS})^n}{n!} e^{-\lambda_{p-1}^+ x_p^{RS}}, \quad n = 0, 1, 2, \dots \quad (5.3)$$

Removing the condition on N from (5.2) by using (5.3), we get

$$E[e^{-sR_p^{RS}} | x_p^{RS}] = e^{-(s + \lambda_{p-1}^+ - \lambda_{p-1}^+ \Theta_{p-1}^+(s)) x_p^{RS}} \quad (5.4)$$

We next consider the period R_p^{RP} . Clearly, the period R_p^{RP} equals C_p^{RP} . It follows from (3.8) and (5.4) that, if the RP phase for a packet of class p is served according to the preemptive repeat-identical discipline, we obtain the LST $R_p^*(s)$ for the residence time R_p as

$$R_p^*(s) = \int_0^\infty \int_0^\infty \int_0^\infty \frac{(s + \lambda_{p-1}^+) e^{-(s + \lambda_{p-1}^+) x} e^{-(s + \lambda_{p-1}^+ - \lambda_{p-1}^+ \Theta_{p-1}^+(s)) y} e^{-sz}}{s + \lambda_{p-1}^+ - \lambda_{p-1}^+ \Theta_{p-1}^+(s) [1 - e^{-(s + \lambda_{p-1}^+) x}]} dB_p^{RP, RS, NP}(x, y, z) \quad (5.5)$$

where $B_p^{RP,RS,NP}(x,y,z)$ is the joint DF for x_p^{RP} , x_p^{RS} , and x_p^{NP} .

Replacing $\Theta_{p-1}^+(s)$ by 1 in (5.5), we can obtain the LST $\bar{x}_p^*(s)$ for the gross service time \bar{x}_p as

$$\bar{x}_p^*(s) = \int_0^\infty \int_0^\infty \int_0^\infty \frac{(s + \lambda_{p-1}^+) e^{-(s + \lambda_{p-1}^+)x}}{s + \lambda_{p-1}^+ e^{-(s + \lambda_{p-1}^+)x}} e^{-s(y+z)} dB_p^{RP,RS,NP}(x,y,z) \quad (5.6)$$

From (5.6) we have

$$E[\bar{x}_p] = \frac{1}{\lambda_{p-1}^+} \left(E[e^{\lambda_{p-1}^+ x_p^{RP}}] - 1 \right) + E[x_p^{RS+NP}] \quad (5.7)$$

$$\begin{aligned} E[\bar{x}_p^2] &= \frac{2}{(\lambda_{p-1}^+)^2} \left\{ E[(e^{\lambda_{p-1}^+ x_p^{RP}})^2] - E[e^{\lambda_{p-1}^+ x_p^{RP}}] - \lambda_{p-1}^+ E[x_p^{RP} e^{\lambda_{p-1}^+ x_p^{RP}}] \right\} \\ &\quad + E[(x_p^{RS+NP})^2] + \frac{2}{\lambda_{p-1}^+} \left(E[x_p^{RS+NP} e^{\lambda_{p-1}^+ x_p^{RP}}] - E[x_p^{RS+NP}] \right) \end{aligned} \quad (5.8)$$

where $x_p^{RS+NP} \equiv x_p^{RS} + x_p^{NP}$. From (5.5), we also get

$$E[R_p] = \frac{E[\bar{x}_p] - \rho_{p-1}^+ E[x_p^{NP}]}{1 - \rho_{p-1}^+} \quad (5.9)$$

6. Waiting Time

We proceed to find the distribution of the waiting time W_p , which is defined as the interval from the moment of arrival of a packet of class p to the moment at which it first receives service. Note that the completion time C_p also represents the time from the start of service to a packet of class p until the first moment when another packet of class p can enter the server. Therefore, the waiting time of a packet of class p that arrives during the delay cycle initiated with D is identical to the waiting time of a packet that arrives during the delay cycle in a nonpriority M/G/1 system with arrival rate λ_p and service time distribution $C_p^*(s)$. Thus we have

$$W_p^*(s|D\text{-delay cycle}) = \frac{(1 - \lambda_p E[C_p]) (1 - D^*(s))}{E[D] (s - \lambda_p + \lambda_p C_p^*(s))} \quad (6.1)$$

where $D^*(s)$ is the LST of the DF for D [5, Section 8.4] [8, Section 5.10]. Let us assume that at a random point in time, the system is in a D -delay cycle with probability P_D , or in an idle period with probability P_0 so that P_D and P_0 satisfy the following condition:

$$P_0 + \sum_D P_D = 1 \quad (6.2)$$

Because Poisson arrivals see time averages (PASTA) [12], the unconditioned waiting time is given by

$$W_p^*(s) = P_0 W_p^*(s|\text{idle}) + \sum_D P_D W_p^*(s|D\text{-delay cycle}) \quad (6.3)$$

We will obtain $W_p^*(s)$ by using (6.3). In our priority system, a packet of class p arrives during an idle period with probability $P_0 = 1 - \rho$ and its waiting time is zero. Therefore the LST of the DF for the waiting time of this packet is given by

$$W_p^*(s|\text{idle}) = 1; \quad P_0 = 1 - \rho \quad (6.4)$$

A packet of class p arrives during a $B_p^{-,RP+RS}$ -period, which is the RP or RS phase of the service time for a packet of class $p+1, \dots, P$, with probability $P_p^{-,RP+RS} = \rho - \rho_p^+ - \rho_p^{-,NP}$, where $\rho_p^{-,NP} = \sum_{k=p+1}^P \lambda_k \int_0^\infty x dB_k^{NP}(x)$ and $B_k^{NP}(x)$ is the DF for x_k^{NP} . Obviously, its waiting time is again zero. Therefore we have

$$W_p^*(s|B_p^{-,RP+RS}\text{-period}) = 1; \quad P_p^{-,RP+RS} = \rho - \rho_p^+ - \rho_p^{-,NP} \quad (6.5)$$

In addition, there are three types of delay cycles in the system: the delay cycle starting with a B_{p-1}^+ -cycle, B_p -cycle, and $B_p^{-,NP}$ -cycle. The B_{p-1}^+ -cycle is a delay cycle initiated with the service time for a packet of class $1, 2, \dots, p-1$, the B_p -cycle is a delay cycle initiated with the service time for a packet of class p , and $B_p^{-,NP}$ -cycle is a delay cycle initiated with the NP phase of the service time for a packet of class $p+1, \dots, P$. Let P_{p-1}^+, P_p and $P_p^{-,NP}$ be the probabilities that a packet of class p arrives during a B_{p-1}^+ -cycle, B_p -cycle, and $B_p^{-,NP}$ -cycle, respectively. They should satisfy the relation

$$P_{p-1}^+ + P_p + P_p^{-,NP} = \rho_p^+ + \rho_p^{-,NP} \quad (6.6)$$

The probabilities P_{p-1}^+, P_p , and $P_p^{-,NP}$ can be found as follows. First, note that the mean number of times that the system enters the $B_p^{-,NP}$ -cycle per unit time is λ_p^- . The mean length of a $B_p^{-,NP}$ -cycle is $E[x_p^{-,NP}]/(1 - \rho_p^+)$, where $E[x_p^{-,NP}] = \rho_p^{-,NP}/\lambda_p^-$. Thus

$$P_p^{-,NP} = \lambda_p^- \times \frac{E[x_p^{-,NP}]}{1 - \rho_p^+} = \frac{\rho_p^{-,NP}}{1 - \rho_p^+} \quad (6.7)$$

From (6.6) and (6.7), we get

$$P_{p-1}^+ + P_p = \frac{\rho_p^+(1 - \rho_p^+ - \rho_p^{-,NP})}{1 - \rho_p^+} \quad (6.8)$$

The system enters the B_{p-1}^+ -cycle or the B_p -cycle only from the state in which there are no packets of class $1, 2, \dots, p-1$ in the system. Therefore, the ratio of P_{p-1}^+ to P_p equals

$$\frac{P_{p-1}^+}{P_p} = \frac{\rho_{p-1}^+}{\rho_p} \quad (6.9)$$

From (6.8) and (6.9), we determine

$$P_{p-1}^+ = \frac{\rho_{p-1}^+(1 - \rho_p^+ - \rho_p^{-,NP})}{1 - \rho_p^+} \quad (6.10)$$

$$P_p = \frac{\rho_p(1 - \rho_p^+ - \rho_p^{-,NP})}{1 - \rho_p^+} \quad (6.11)$$

The LST of the DF for the waiting time of a packet of class p that arrives during the $B_p^{-,NP}$ -cycle is obtained by substituting $D^*(s) = B_p^{*(-,NP)}[s + \lambda_{p-1}^+ - \lambda_{p-1}^+ \Theta_{p-1}^+(s)]$ into (6.1), where $B_p^{*(-,NP)}(s) = \frac{1}{\lambda_p^-} \sum_{k=p+1}^P \lambda_k \int_0^\infty e^{-sx} dB_k^{NP}(x)$. Thus we get

$$W_p^*(s|B_p^{-,NP}\text{-cycle}) = \frac{\lambda_p^-(1 - \rho_p^+) \left[1 - B_p^{*(-,NP)}(\sigma_{p-1}) \right]}{\rho_p^{-,NP} \left[s - \lambda_p + \lambda_p C_p^*(s) \right]} \quad (6.12)$$

where $\sigma_{p-1} \equiv s + \lambda_{p-1}^+ - \lambda_{p-1}^+ \Theta_{p-1}^+(s)$. Similarly, the LST of the DF for the waiting time of a packet of class p that arrives during the B_{p-1}^+ -cycle and B_p -cycle are respectively given by

$$W_p^*(s|B_{p-1}^+\text{-cycle}) = \frac{\lambda_{p-1}^+(1 - \rho_p^+) (1 - \Theta_{p-1}^+(s))}{\rho_{p-1}^+ [s - \lambda_p + \lambda_p C_p^*(s)]} \quad (6.13)$$

$$W_p^*(s|B_p\text{-cycle}) = \frac{\lambda_p(1 - \rho_p^+) (1 - C_p^*(s))}{\rho_p [s - \lambda_p + \lambda_p C_p^*(s)]} \quad (6.14)$$

We now apply (6.3) to obtain $W_p^*(s)$ as follows:

$$W_p^*(s) = \frac{(1 - \rho_p^+ - \rho_p^{-,NP}) [s + \lambda_{p-1}^+ - \lambda_{p-1}^+ \Theta_{p-1}^+(s)] + \lambda_p^- [1 - B_p^{*(-,NP)}(\sigma_{p-1})]}{s - \lambda_p + \lambda_p C_p^*(s)} \quad (6.15)$$

Then we obtain the mean waiting time as

$$E[W_p] = \frac{\lambda_p^+(1 - \rho_p^+)^2 E[(\Theta_p^+)^2]}{2(1 - \rho_{p-1}^+)} + \frac{\lambda_p^- E[(x_p^{-,NP})^2]}{2(1 - \rho_{p-1}^+)(1 - \rho_p^+)} \quad (6.16)$$

where we have used (4.4).

7. Response Time and Queue Size

The response time T_p of a packet of class p consists of two independent random variables, the waiting time W_p and the residence time R_p . Thus the LST of the DF for the T_p given by

$$T_p^*(s) = W_p^*(s)R_p^*(s) \quad (7.1)$$

From (7.1), we have

$$E[T_p] = \frac{\lambda_p^+(1 - \rho_p^+)^2 E[(\Theta_p^+)^2]}{2(1 - \rho_{p-1}^+)} + \frac{\lambda_p^- E[(x_p^{-,NP})^2]}{2(1 - \rho_{p-1}^+)(1 - \rho_p^+)} + \frac{E[\bar{x}_p] - \rho_{p-1}^+ E[x_p^{NP}]}{1 - \rho_{p-1}^+} \quad (7.2)$$

where $E[W_p]$ and $E[R_p]$ have been given in Sections 3 and 4 for each preemptive repeat discipline.

Let $\Pi_p(z)$ be the probability generating function of the number of packets of class p in the system at the departure time of a class p . We then have the following relationship

$$\Pi_p(z) = W_p^*(\lambda_p - \lambda_p z)R_p^*(\lambda_p - \lambda_p z) \quad (7.3)$$

which represents the number of packets of class p that arrived while the arbitrary departing packet of class p was in the system. Note that $\Pi_p(z)$ is also the probability generating function of the number L_p of packets of class p in the system at an arbitrary time. This comes from PASTA (Poisson arrivals see time averages) and Burke's theorem on the process with unit jumps applied to the number of packets of class p present in the system [6, Section 5]. Then the mean queue size $E[L_p]$ of packets of class p at an arbitrary time as

$$E[L_p] = \frac{\lambda_p \lambda_p^+(1 - \rho_p^+)^2 E[(\Theta_p^+)^2]}{2(1 - \rho_{p-1}^+)} + \frac{\lambda_p \lambda_p^- E[(x_p^{-,NP})^2]}{2(1 - \rho_{p-1}^+)(1 - \rho_p^+)} + \frac{\lambda_p (E[\bar{x}_p] - \rho_{p-1}^+ E[x_p^{NP}])}{1 - \rho_{p-1}^+} \quad (7.4)$$

8. Numerical Examples

To investigate numerically the performance of the system, we consider a system with 5 classes of packets. We assume that all λ'_p 's are identical and that their RP and NP phases of the service times for each class are 10^{-3} seconds (constant) respectively. This corresponds to the assumption that the lengths of the header and the trailer of packets are fixed in the application example mentioned in Section 1. In the first example, we assume that the RS phase of the service time (the information field of a packet) of each class is also constant, which is 10^{-2} second. In the second example, we assume that the RS phase of the service time of each class is exponentially distributed with rates $10^2/\text{sec}$ (the mean is 10^{-2} seconds). The mean response times for each class of packet have been computed, and are shown in Figures 1 and 2 against the total server utilizations ρ .

As evident in these figures, if we compare the response times for the two different DFs of the RS phase of service time, the response times for each class of packets with exponentially distributed RS phases are greater and more discriminative than those with constant RS phases.

Acknowledgment

This research is supported in part by University of Tsukuba Research Project Foundation.

References

- [1] J.D. Atkins, Path control-The network layer of System Network Architecture, In: *Computer Network Architectures and Protocols*, edited by P.E. Green, Jr., (Plenum Press, New York, 1982) 297-326.
- [2] B. Avi-Itzhak, I. Brosh, and P. Naor, On discretionary priority queueing, *Zeitschrift fuer Angewandte Mathematik und Mechanik*, Vol.6 (1964) 235-318.
- [3] Y.Z. Cho, An efficient priority-scheduling algorithm for integrated services packet networks, Ph.D. dissertation, Korea Advanced Institute of Science and Technology, 1988.
- [4] Y.Z. Cho and C.K. Un, Analysis of the M/G/1 queue under a combined preemptive/nonpreemptive priority discipline, *IEEE Transactions on Communications*, Vol.41, No.1 (1993) 132-141.
- [5] R.M. Conway, W.L. Maxwell, and L.W. Miller, *Theory of Scheduling*, (Addison-Wesley, Reading, Massachusetts, 1967).
- [6] R.B. Cooper, *Introduction to Queueing Theory*, (Third edition, CEEPress Books, Washington D.C., 1990).
- [7] N.K. Jaiswal, *Priority Queues*, (Academic Press, New York, 1968).
- [8] L. Kleinrock, *Queueing Systems, Volume 1: Theory*, (John Wiley and Sons, New York, 1975).
- [9] M. Komatsu, A Priority queue with two priority classes of which non-priority class call has preemptive and nonpreemptive parts,(in Japanese) *The Transactions of IEICE*, Vol. J71-A, No. 11 (1988) 2027-2032, (Correction in Vol. J72-A, No. 3 (1989) p. 617).
- [10] M. Stallings, *Local and Metropolitan Area Networks*, (Fourth edition, Macmillan, New York, 1983).
- [11] H. Takagi, *Queueing Analysis, A Foundation of Performance Evaluation, Volume 1: Vacation and Priority Systems*, (Elsevier, Amsterdam, 1991).
- [12] R.W. Wolff, Poisson arrivals see time averages, *Operations Research*, 30 (1982) 223-231.

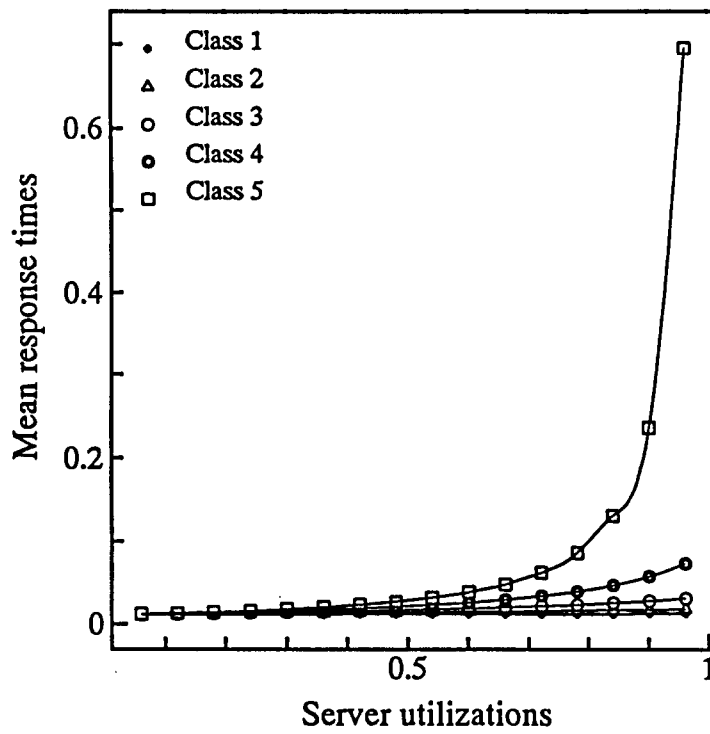


Figure 3: RS phase : 0.01s(constant)
 RP and NP phases : 0.001s(constant)

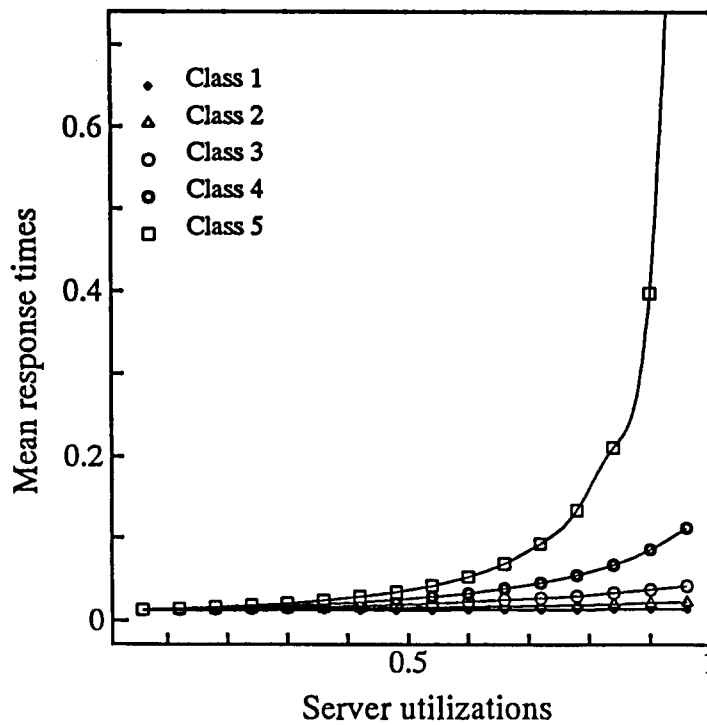


Figure 4: RP and NP phases : 0.001s(constant)
 RS phase : exponential dis. with rate 100/s